

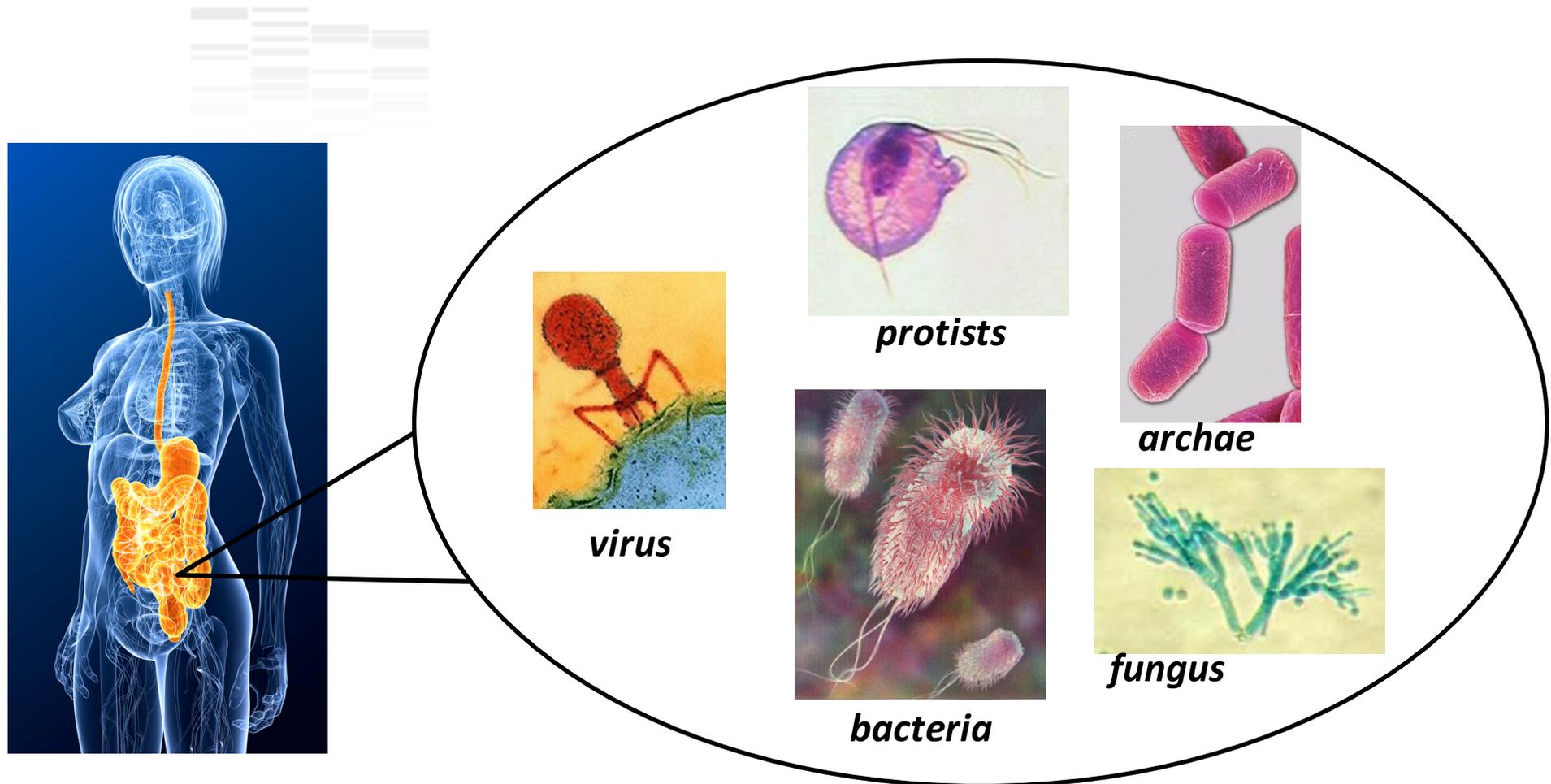


Le langage Go dans le calcul HPC Big Data

Louiza HANIS
INRA, stagiaire M2 master Haute Performance & Simulation



La métagénomique ?

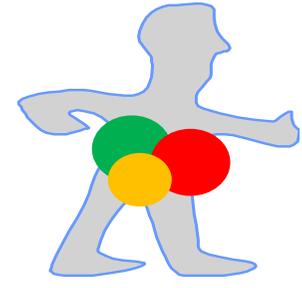
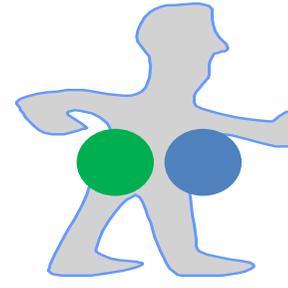
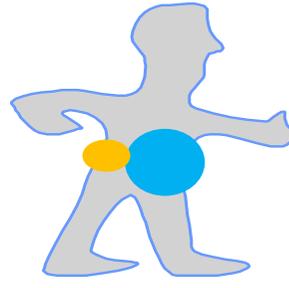
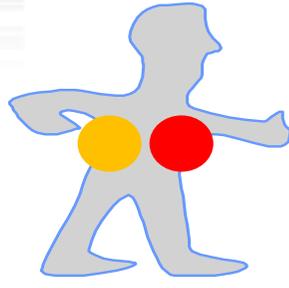


- **Big Data de l'ADN Bactérien**
- **2 kg – plus de bactéries que de cellules humaines**
- **Comment découvrir la ou les bactéries signatures (Biomarqueurs)?**

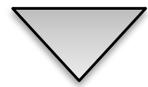
Métagénomique quantitative



Flore intestinale



Échantillons de
fèces



Analyse globale de
l'ADN



Identification & Quantification



Abondance relative par
projection sur un catalogue
de références

1
0
1
1
2
0
0
0
0
....

0
0
5
0
0
8
0
0
0

0
10
0
0
0
0
10
1

1
0
20
1
1
0
0
1

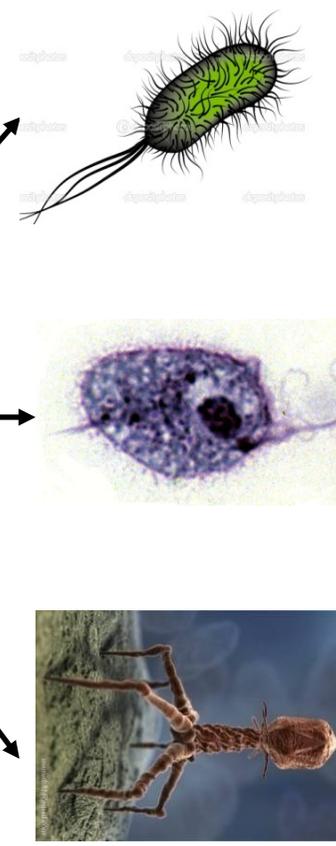
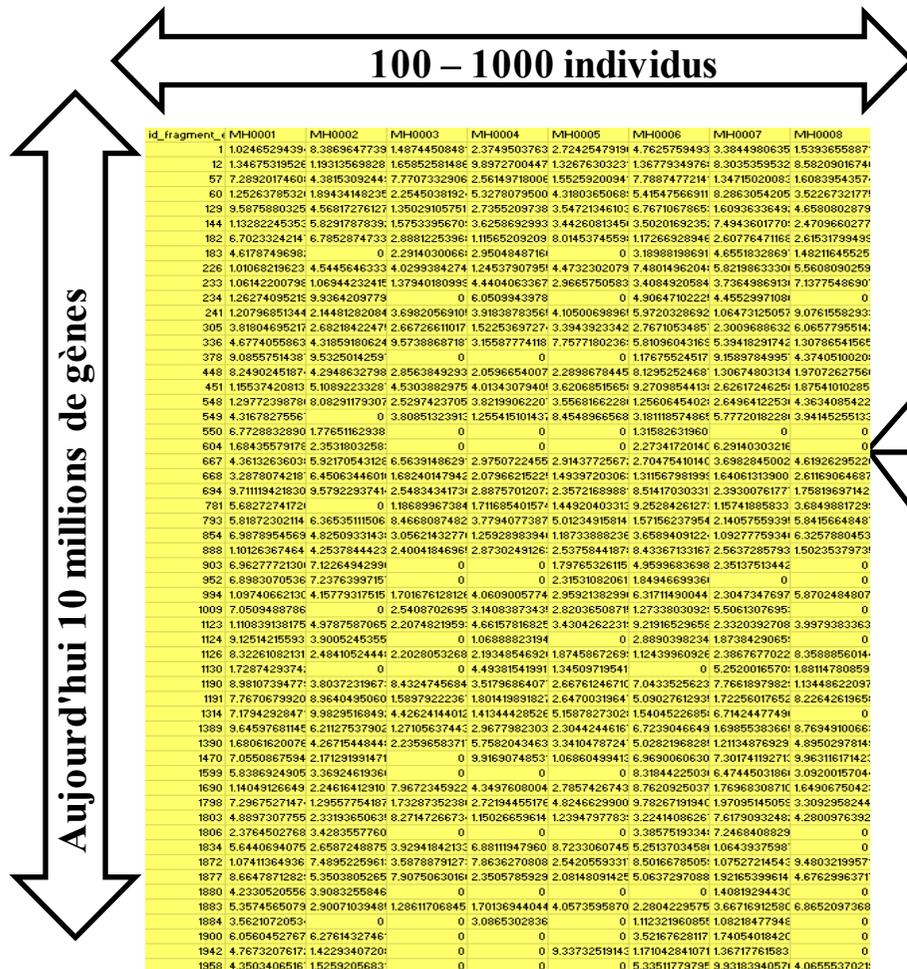
Il s'agit d'extraire des clusters de manière déterministe et exhaustive



bactérie = ~3000 gènes

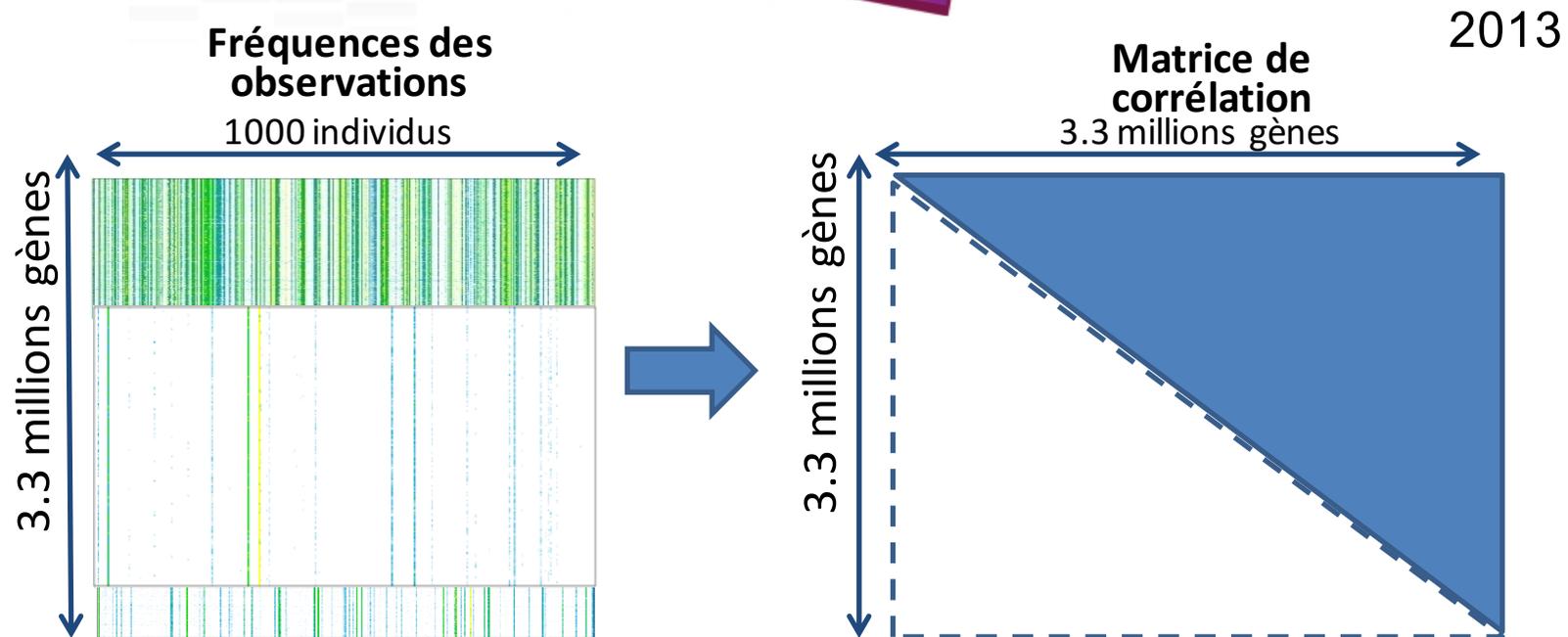
parasites = ~6000 gènes

Virus = ~50 gènes



Maintenir un code HPC spécifique?

- IN/OUT



EN 2010 un code de calcul OpenCL/CUDA a été réalisé pour traiter des matrices de taille 1000 * 3300000

http://www.genci.fr/sites/default/files/INRA_AS_Plus_Tello_Almeida_Metaprof.pdf

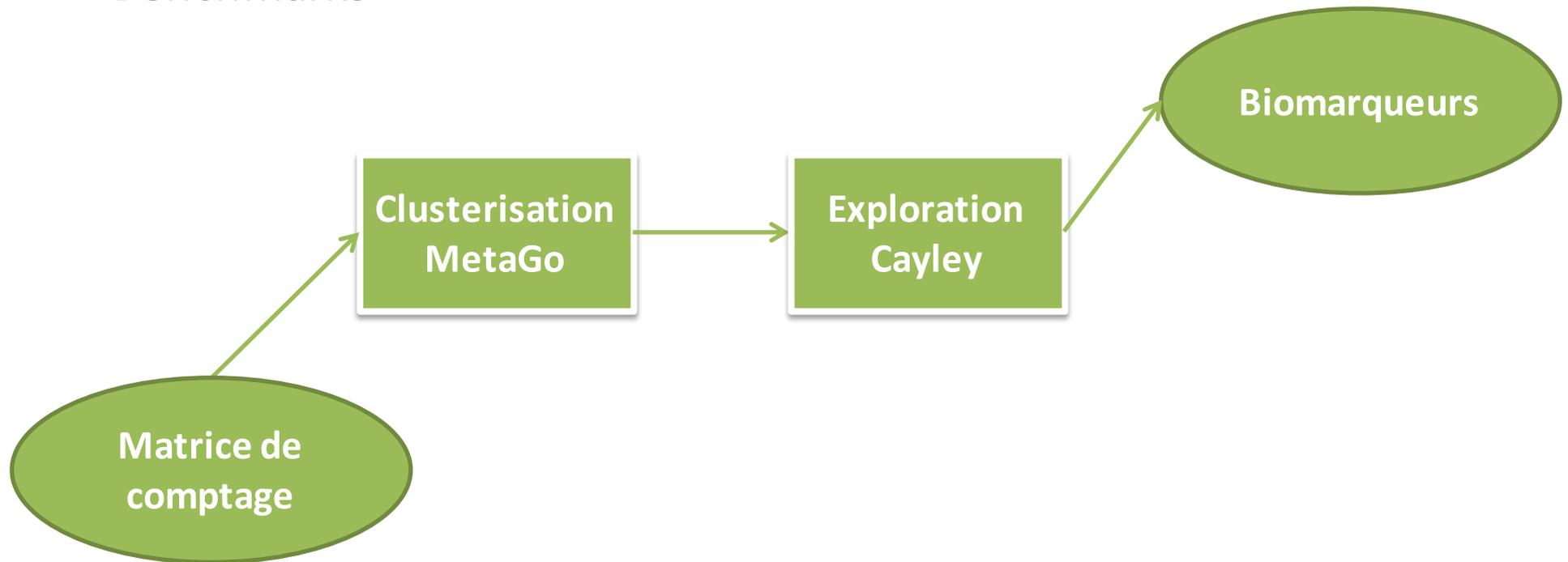
→ Comment maintenir du code HPC spécifique BigData (CUDA) et du code de traitement de data en amont de R?

→ Peut-on unifier les codes?

Plan de développement



- ✓ Tests unitaires
- ✓ Tests d'intégration
- ✓ Tests de vérification avec R
- ✓ Benchmarks



Les contraintes

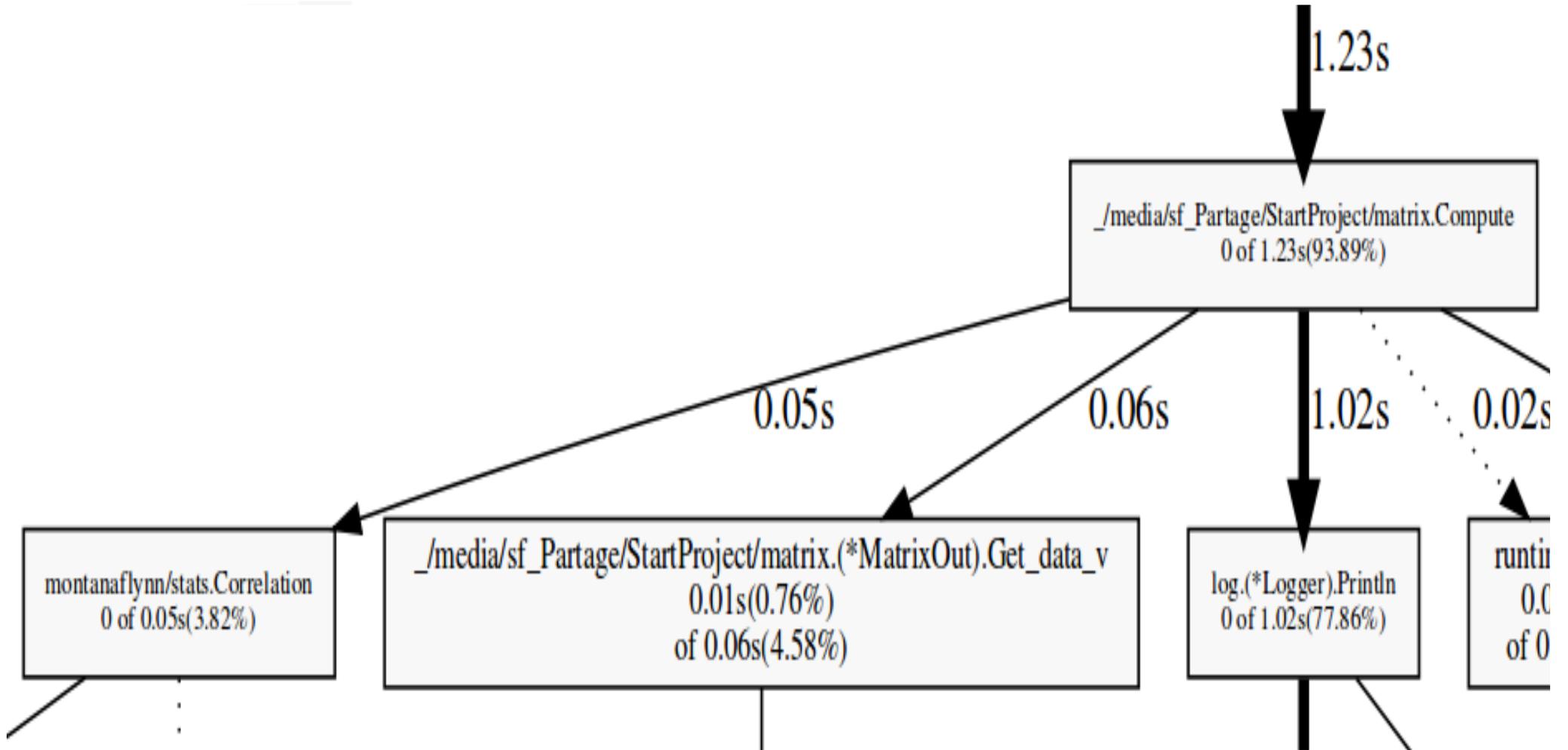


Langage R

- Langage interprété \equiv Langage lent
- Pas orthogonal
- Non conçu pour le Big Data

Go Profiler

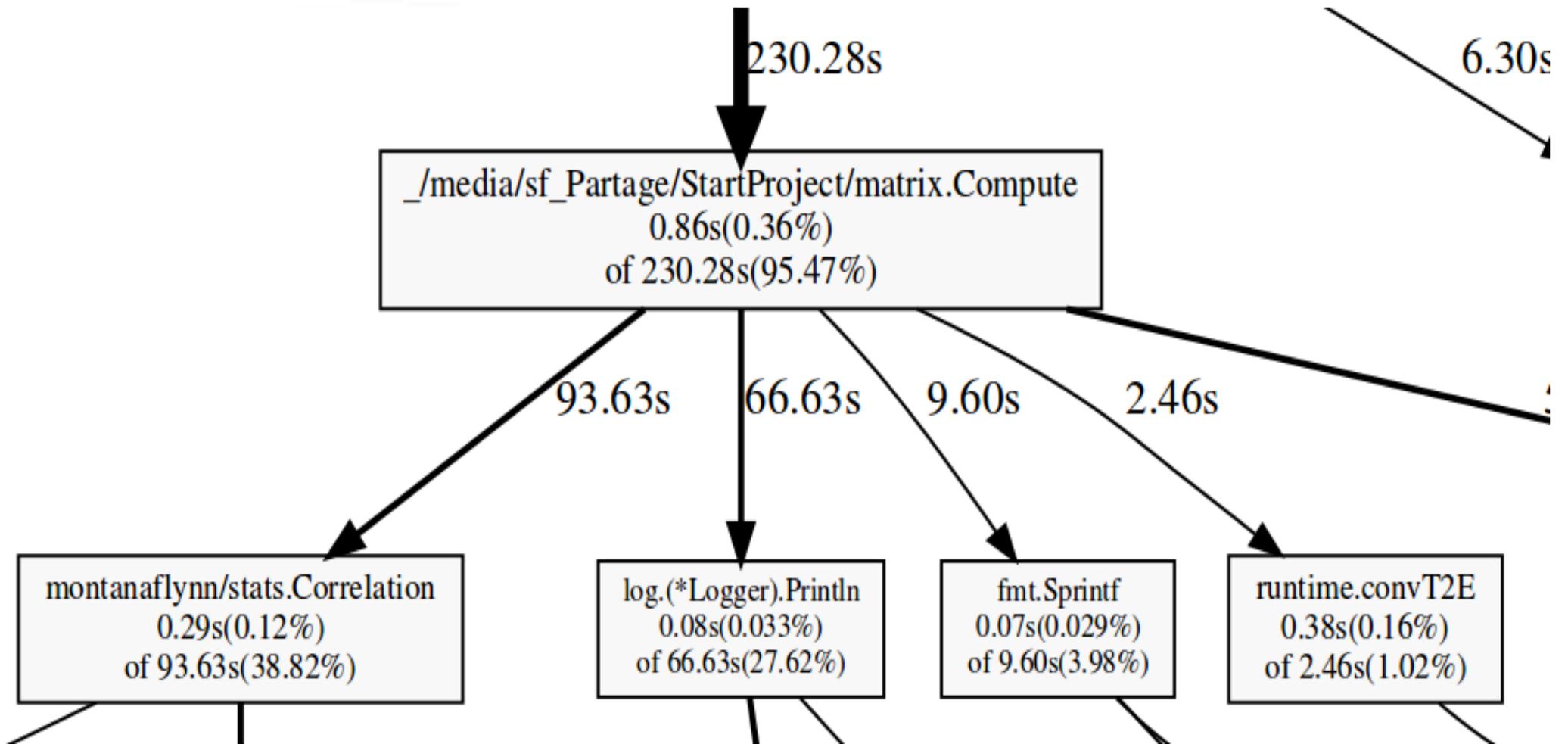
❖ Temps d'écriture > Temps de calcul



Matrice 1000 * 50

Go Profiler

❖ Temps de calcul > Temps d'écriture



Matrice 5000 * 1267

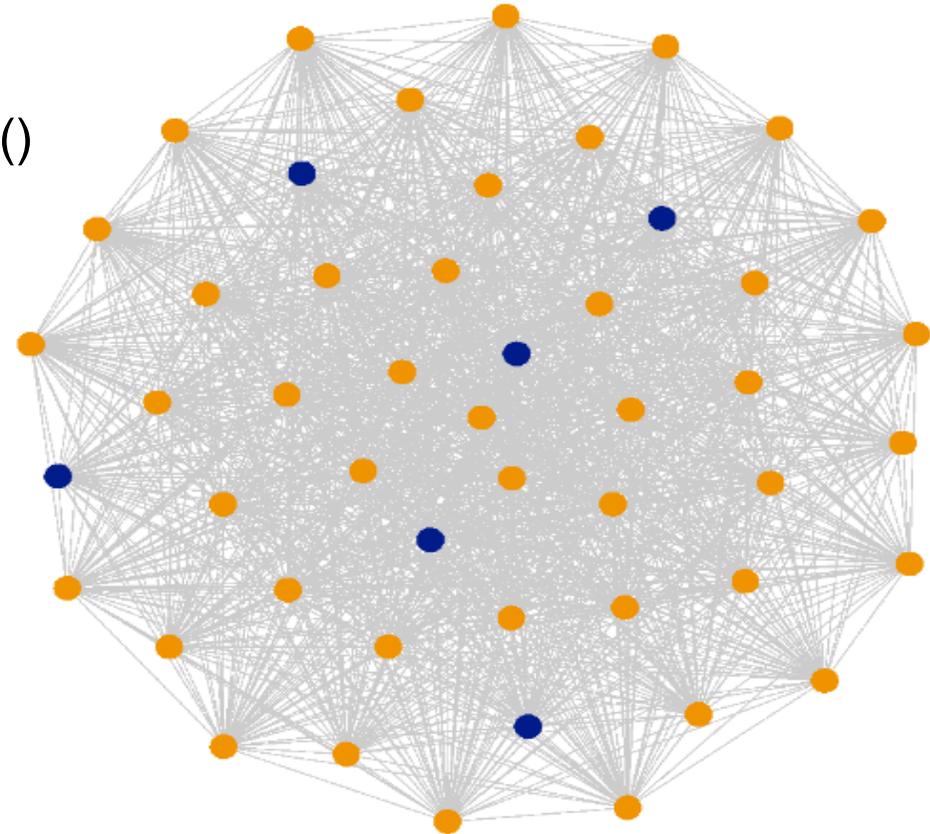
Cayley

- <https://github.com/google/cayley>
- Il s'agit d'un navigateur (http) et d'un moteur de requête pour une database de type "graphe orienté"
- Stockage par triplet (sujet, prédicat, objet)
 - Sujet
 - Prédicat/Objet = règle+objet (ex. Prolog)
- Un ensemble de triplet = RDF Graph (Resource Description Framework)
- Triplet+"label" = N-quads
- Repose sur Bolt (Clé/Valeur + B-tree)

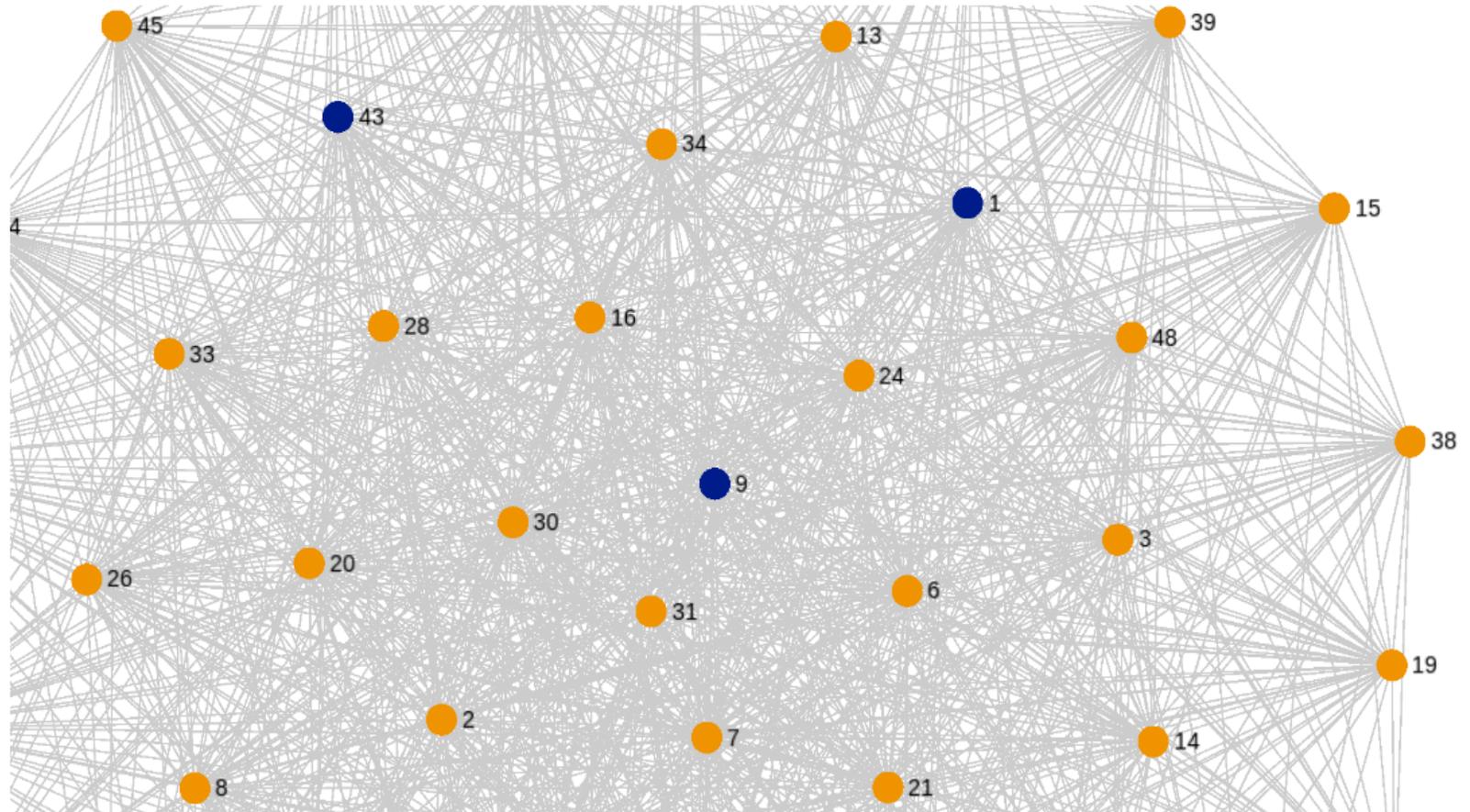
Cayley



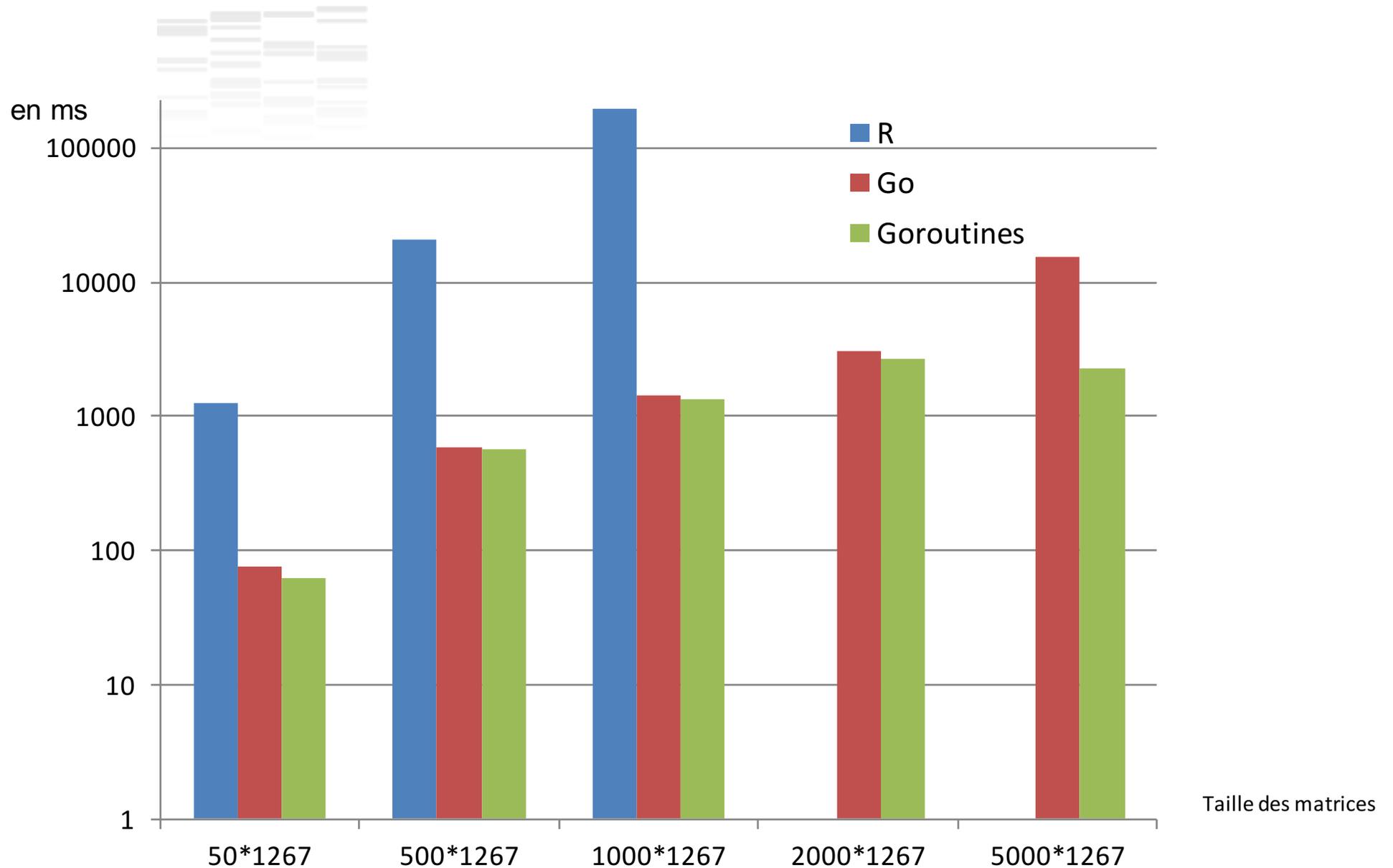
```
g.V().Tag("source").Out().Tag("target").All()
```



Cayley



Le benchmark R vs Go



Avantages de développer en go



- ✓ Environnement complet
- ✓ Intégration de plusieurs outils
- ✓ Packages faciles à importer (Go get...)
- ✓ Documentation riche